

# BIG DATA ANALYSIS USING SHRINKAGE STRATEGIES

BAHADIR YÜZBAŞI

*Department of Econometrics, Inonu University, Turkey*

MOHAMMAD ARASHI

*Department Statistics, Shahrood University of Technology, Shahrood, Iran*

S. EJAZ AHMED

*Department of Mathematics and Statistics, Brock University, Canada*

**ABSTRACT.** In this paper, we apply shrinkage strategies to estimate regression coefficients efficiently for the high-dimensional multiple regression model, where the number of samples is smaller than the number of predictors. We assume in the sparse linear model some of the predictors have very weak influence on the response of interest. We propose to shrink estimators more than usual. Specifically, we use integrated estimation strategies in sub and full models and shrink the integrated estimators by incorporating a bounded measurable function of some weights. The exhibited double shrunken estimators improve the prediction performance of sub models significantly selected from existing Lasso-type variable selection methods. Monte Carlo simulation studies as well as real examples of eye data and Riboavin data confirm the superior performance of the estimators in the high-dimensional regression model.

**Keywords:** Double shrinkage; High-dimension; Penalty; Prediction; Sparse regression model.

## 1. INTRODUCTION

Nowadays many researchers have focused on the analysis of big data, because of existence trend in computer science and statistics. Comparing to the usual datasets, big data refer to high-dimensional, unusual and unstructured data. The analysis of big data needs methods other than traditional analytical frameworks.

Let  $n$  denote the sample size or number of observations and  $p$  the number of features or variables. Data scientists consider big data as when  $n$  is too large, however in medical and genetic researches, engineering and financial studies, one mostly involves small  $n$  large  $p$  problem, known as high-dimensional data. As Pyne et al. (2016) pointed in big data analytics, some domains of big data such as finance or health do even produce infinite dimensional functional data, which are observed not as points but functions, such as growth curves, online auction bidding trends, etc. As Wang et al. (2012) pointed, big data are data on a massive scale in terms of volume, intensity, and complexity that exceed the capacity of standard analytic tools. Ahmed (2014b) collected some research

---

*E-mail addresses:* (Bahadır Yüzbaşı) b.yzb@hotmail.com, (Mohammad Arashi) m.arashi\_stat@yahoo.com, (S. Ejaz Ahmed) sahmed5@brocku.ca.

contributions in the field of big data analytics to highlight high-dimensional methods in big data challenges.

As [Ahmed and Yüzbaşı \(2016a\)](#) and [Ahmed and Yüzbaşı \(2016b\)](#) pointed, the term “big data” is not well defined, but its problems are real and statisticians need to play a more important role in this arena. The big data or data science is an emerging field. In 2013, American Statistical Association (ASA) proposed three reasons to show ASA has not been very involved in big data. President of Institute of Mathematical Statistics (IMS), Bin Yu, in 2014 called for statisticians to own data science by working on real problems such as those from genomics, neuroscience, astronomy, nanoscience, computational social science, personalized medicine/healthcare, finance, and government; relevant methodology/theory will follow naturally.

There is an increasing demand for efficient prediction strategies for analyzing high-dimensional data in big data streams. For example, data arising from gene expression arrays, social network modeling, clinical, genetics and phenotypic data. Due to the trade-off between model complexity and model prediction, the statistical inference of model selection becomes an extremely important and challenging problem in high-dimensional data analysis. Over the past two decades, many penalized regularization approaches have been developed to do variable selection and estimation simultaneously. Among them, least absolute shrinkage and selection operator (LASSO) is one of the recent approaches, [Tibshirani \(1996\)](#). It is a useful technique due to its convexity and computation efficiency. The LASSO is based on squared error and a penalty proportional to regression parameters. [Schellldorfer et al. \(2011\)](#) provides a comprehensive summary of the consistency properties of the LASSO. [Efron et al. \(2004\)](#) introduced the least angle regression algorithm which is a very fast way to draw the entire regularization path for a LASSO estimate of the regression parameters. The penalized likelihood methods have been extensively studied in the literature, see for example, [Tran \(2011\)](#), [Huang et al. \(2008\)](#), [Kim et al. \(2008\)](#), [Wang and Leng \(2012\)](#), [Yuan and Lin \(2006\)](#), [Leng et al. \(2006\)](#), and [Tibshirani et al. \(2005\)](#). The penalized likelihood methods have a close connection to Bayesian procedures. Thus the LASSO estimate corresponds to a Bayes method that puts a Laplacian (double exponential) prior on the regression coefficients. Recent results ([Armagan et al. \(2013\)](#), [Bhattacharya et al. \(2012\)](#), and [Carvalho et al. \(2010\)](#)) have demonstrated that better desirable results can be obtained by using priors with heavier tails than the double exponential prior, in particular, priors with polynomial tails. Our study has concentrated on the widely recognized penalty estimators LASSO and adaptive LASSO (ALASSO). Very recently, [Yüzbaşı and Arashi \(2016\)](#) have proposed double shrinking concept to improve the prediction accuracy of LASSO. Here, we specifically implement the double shrunken estimator on ALASSO.

Following [Ahmed and Yüzbaşı \(2016a\)](#), we consider the estimation problem of regression parameters when there are many potential predictors in the initial/working model and:

- (1) most of them may not have any influence (sparse signals) on the response of interest
- (2) some of the predictors may have strong influence (strong signals) on the response of interest
- (3) some of them may have weak-moderate influence (weak-moderate signals) on the response of interest

It is possible that there may be extraneous predictors in the model. Suppose if the main concern is treatment effect, or the effect of biomarkers, extraneous nuisance variables may be lab effects when several labs are involved, or the age and sex of patients. The analysis will be more precise if “nuisance variables” can be left out of the model. This leads to the consideration of two models: the full model that includes all predictors and possible

extraneous variables, and a candidate submodel that includes the predictors of main concern while leaving out extraneous variables. Further, it is important that we do not automatically remove all the predictors with weak signals from the model. This may result in selecting a biased submodel. A logical way to deal with this framework is to use pretest model selection and estimation strategies that test whether the coefficients of the extraneous variables are zero and then estimate parameters in the model that include coefficients that are rejected by the test. Another strategy is to use Stein-type shrinkage estimators where the estimated regression coefficient vector is shrunk in the direction of the candidate subspace. This “soft threshold” modification of the pretest method has been shown to be efficient in various frameworks. [Ahmed et al. \(2012\)](#), among others have investigated the properties of shrinkage and pretest methodologies for host models.

The model and some estimators are introduced in Section 2. In Section 3 we show-case our suggested estimation strategy. The results of a simulation study that includes comparison of suggested estimator with the penalty estimators are reported in Section 4. Application to real data sets is given in Section 5. Finally, we offer concluding remarks in Section 6.

## 2. ESTIMATION STRATEGIES

In this communication, we consider a high-dimensional linear regression sparse model:

$$(2.1) \quad y_i = \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i, \quad 1 \leq i \leq n \ll p$$

where  $y_i$  observed response variable with predictors  $x_i$ s, and  $\beta_j$  are the regression parameters. Further,  $\varepsilon_i$ s are independent and identically distributed random errors with center 0 and variance  $\sigma^2$ . Similar to most of LASSO penalty-type models, in our approach, we assume the true model is sparse in the sense that most of regression coefficients are zeros except for a few ones and all nonzero  $\beta_j$ 's are larger than noise level,  $c\sigma\sqrt{(2/n)\log(d)}$  with  $c \geq 1/2$ . We refer to [Zhao and Yu \(2006\)](#), [Huang et al. \(2008\)](#), and [Bickel et al. \(2009\)](#) for some insights. In general, the LASSO penalty turns to select an over-fitted model since it penalizes all coefficients equally ([Leng et al. \(2006\)](#)). In reviewed literature several modification and methodologies have been suggested to improve the prediction accuracy for LASSO strategy. For example, the SCAD ([Fan and Li \(2001\)](#)), adaptive LASSO, ([Zou \(2006\)](#)), MCP ([Zhang \(2010\)](#)) and Stein-type LASSO ([Yüzbaşı and Arashi \(2016\)](#)) and several others. These methods select a submodel by shrinking some regression coefficients to zero and provide shrinkage estimators of the remaining coefficients. However, these methods may force the relatively more weak coefficients towards zeros as compared to LASSO, resulting in under-fitted models subject to a much larger selection bias in the presence of significant number of weak signals.

Following [Ahmed and Yüzbaşı \(2016a\)](#) and [Ahmed and Yüzbaşı \(2016b\)](#), in this paper, we consider the estimation and prediction problem for the sparse regression models when there are many potential predictors that have weak influence on the response of interest. The analysis will be relatively more precise if “weak effect” variables can be weighted for the ultimate model prediction. This leads to the consideration of two models: the over-fitted model that includes predictors with strong signals and possibly some predictors with weak signals selected by LASSO. On the other, we select an underfitted model that possibly includes the predictors with strong signals while leaving out predictors with weak effect by using ALASSO. One way to deal with this framework is to use Stein-type shrinkage estimators where the estimated regression coefficient vector is shrunk in the direction of the under-fitted model. This “soft threshold” modification of the

pretest method has been shown to be efficient in various frameworks. [Saleh \(2006\)](#) and [Ahmed et al. \(2012\)](#), among others have examined the properties of Stein-type shrinkage estimation strategies for a host of models.

Consider the following regression model

$$(2.2) \quad \mathbf{Y} = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where  $\mathbf{Y} = (y_1, y_2, \dots, y_n)'$  is a vector of responses,  $\mathbf{X}_n$  is an  $n \times p$  fixed design matrix,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  is an unknown vector of parameters,  $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)'$  is the vector of unobservable random errors, and the superscript  $(\cdot)'$  denotes the transpose of a vector or matrix. We do not make any distributional assumption about the errors except that  $\boldsymbol{\varepsilon}$  has a cumulative distribution function  $F(\boldsymbol{\varepsilon})$  with  $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ , and  $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \sigma^2 \mathbf{I}$ , where  $\sigma^2$  is finite.

For  $n > p$  the classical estimator of  $\boldsymbol{\beta}$  by minimizing the least square function and is given by

$$\hat{\boldsymbol{\beta}}^{\text{LSE}} = (\mathbf{X}_n' \mathbf{X}_n)^{-1} \mathbf{X}_n' \mathbf{Y}.$$

However, since we are dealing with a high-dimensional situation, i.e.  $n < p$  so  $(\mathbf{X}'\mathbf{X})^{-1}$  will not exist and thus no solution. However, one can employ the generalized inverse to revert the problem. In the current set-up we are assuming that the model is sparse so it is desirable to use penalized likelihood method to obtain a meaningful solution as was briefly discussed in our Introduction section. Penalty estimators are a class of estimators in the least penalized squares family of estimators, see [Ahmed \(2014a\)](#). This method involves penalizing the regression coefficients, and shrinking a subset of them to zero. In other words, the penalized procedure produces a submodel and subsequently estimates the submodel parameters. Several penalty estimators have been proposed in the literature for linear and generalized linear models. In this section, we consider the LASSO and the ALASSO. By shrinking some regression coefficients to zero, these methods select parameters and estimation simultaneously. [Frank and Friedman \(1993\)](#) introduced bridge regression, a generalized version of penalty (or absolute penalty type) estimators. For a given penalty function  $\pi(\cdot)$  and regularization parameter  $\lambda$ , the general form can be written as

$$S(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}_n \boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}_n \boldsymbol{\beta}) + \lambda \pi(\boldsymbol{\beta}),$$

where the penalty function is of the form

$$(2.3) \quad \pi(\boldsymbol{\beta}) = \sum_{j=1}^m |\beta_j|^\gamma, \quad \gamma > 0.$$

The penalty function in (2.3) bounds the  $L_\gamma$  norm of the parameters in the given model as  $\sum_{j=1}^m |\beta_j|^\gamma \leq t$ , where  $t$  is the tuning parameter that controls the amount of shrinkage. We see that for  $\gamma = 2$ , we obtain ridge estimates which are obtained by minimizing the penalized residual sum of squares

$$(2.4) \quad \hat{\boldsymbol{\beta}}^{\text{Ridge}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\},$$

where  $\lambda$  is the tuning parameter which controls the amount of shrinkage. [Frank and Friedman \(1993\)](#) did not solve for the bridge regression estimators for any  $\gamma > 0$ . Interestingly, for  $\gamma < 2$ , it shrinks the coefficient towards zero, and depending on the value of  $\lambda$ , it sets some of them to be exactly zero. Thus, the procedure combines variable selection and shrinking of the coefficients of penalized regression. [Gao et al. \(2016\)](#) suggested weighted ridge estimator for high dimensional setting, and investigated the advantages of post selection positive part of shrinkage estimators both theoretically and numerically.

An important member of the penalized least squares family is the  $L_1$  penalized least squares estimator, which is obtained when  $\gamma = 1$ , and is called LASSO.

**2.1. LASSO.** The least absolute shrinkage and selection operator was proposed by Tibshirani (1996), which performs variable selection and parameter estimation simultaneously. LASSO is closely related with ridge regression. LASSO solutions are similarly defined by replacing the squared penalty  $\sum_{j=1}^p \beta_j^2$  in the ridge solution (2.4) with the absolute penalty  $\sum_{j=1}^p |\beta_j|$  in the LASSO,

$$(2.5) \quad \hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

Although the change apparently looks subtle, the absolute penalty term made it impossible to have an analytic solution for the LASSO. Originally, LASSO solutions were obtained via quadratic programming. Later, Efron et al. (2004) proposed Least Angle Regression (LAR), a type of stepwise regression, with which the LASSO estimates can be obtained at the same computational cost as that of an ordinary least squares estimation. Further, the LASSO estimator remains numerically feasible for dimensions of  $p$  that are much higher than the sample size  $n$ .

**2.2. Adaptive LASSO.** Zou (2006) modified the LASSO penalty by using adaptive weights on  $L_1$  penalties on the regression coefficients. Such a modified method was referred to as ALASSO. It has been shown theoretically that the ALASSO estimator is able to identify the true model consistently, and the resulting estimator is as efficient as the oracle.

The ALASSO of  $\hat{\beta}^{\text{ALASSO}}$  are obtained by

$$(2.6) \quad \hat{\beta}^{\text{ALASSO}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \hat{w}_j |\beta_j| \right\},$$

where the weight function is

$$\hat{w}_j = \frac{1}{|\hat{\beta}_j^*|^\gamma}; \quad \gamma > 0,$$

and  $\hat{\beta}_j^*$  is a root- $n$  consistent estimator of  $\beta$ . Equation (2.6) is a “convex optimization problem and its global minimizer can be efficiently solved” (Zou, 2006).

The main objective of this research article is to improve the estimation accuracy of the active set of the regression parameters by combining an over-fitted model estimators with an under-fitted one. For this purpose, we follow the methodology of double shrunk estimator of Yüzbaşı and Arashi (2016). As stated earlier, the LASSO produce an over-fitted model as compared with ALASSO and other variable selection methods. The LASSO strategy retains some regression coefficients with weak effects and as well as some with weak effects in the resulted model. On the other hand, aggressive variable selection strategies may force moderate and effects coefficients towards zero, resulting in under-fitted models with a fewer variable of strong effect. The idea here is to combine estimators from an under-fitted model with an over-fitted model using a non-linear shrinkage technique incorporating a measurable bounded function.

### 3. DOUBLE SHRUNKEN ESTIMATORS

In this section, we show how to shrink more the addressed penalized estimators, in the combination of two submodels produced by two distinct variable selection techniques. Similar to Ahmed and Yüzbaşı (2016a), the idea is to work with a sparse model that will be all the predictors included and then apply two variable selection methods with high



and low penalties, respectively. Finally, we combine the estimates from two models to improve post estimation and prediction performances, respectively and incorporate the concept of double shrunken of [Yüzbaşı and Arashi \(2016\)](#).

**3.1. Working Model.** Consider the following high dimensional sparse regression model with strong and weak-to-moderate signals

$$(3.1) \quad \mathbf{Y} = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad p > n$$

Suppose we can divide the index set  $\{1, \dots, p\}$  into three disjoint subsets:  $S_1$ ,  $S_2$  and  $S_3$ . In particular,  $S_1$  includes indexes of nonzero  $\beta_i$ 's which are large and comfortably detectable. The set  $S_2$ , being the intermediate, includes indexes of those nonzero  $\beta_j$  with weak-to-moderate but nonzero effects. By the assumption of sparsity  $S_3$  includes indexes with only zero coefficients and can be easily discarded by exiting variable selection methods. Thus,  $S_1$  and  $S_3$  are able to be retained and discarded by using existing variable selection techniques, respectively. However, it is possible that the  $S_2$  may be covertly included either in  $S_2$  or  $S_3$  depending on existing LASSO-type methods. For the case when  $S_2$  may not be separated from  $S_3$ , some work has been done in this area, see [Zhang and Zhang \(2014\)](#) and others. [Hansen \(2015\)](#) has showed using simulation studies that such a LASSO estimate often performs worse than the post selection least square estimate. To improve the prediction error of a LASSO-type variable selection approach, some (modified) post least squares estimators are studied in [Belloni and Chernozhukov \(2009\)](#) and [Liu and Yu \(2013\)](#).

However, we are interested in cases when covariates in  $S_1$  are kept in the model, and some or all covariates in  $S_2$  are also included in  $S_1$ , which may or may not be useful for prediction purposes. It is possible that one variable selection strategies may produce an over-fitted model, that is retaining predictors from  $S_1$  and  $S_2$ . On the other hand, other methods may produce an under-fitted model keeping only predictors from  $S_1$ . Thus, the predictors in  $S_2$  should be subject to further scrutiny to improve the prediction error.

We partition the design matrix such that  $\mathbf{X} = (\mathbf{X}_{S_1} | \mathbf{X}_{S_2} | \mathbf{X}_{S_3})$ , Further,  $\mathbf{X}_{n1}$  is  $n \times p_1$ ,  $\mathbf{X}_{n2}$  is  $n \times p_2$ , and  $\mathbf{X}_{n3}$  is  $n \times p_3$  submatrix of predictors, respectively; and  $p = p_1 + p_2 + p_3$ . Here we make the usual assumption that  $p_1 \leq p_2 < n$  and  $p_3 > n$ .

Thus, our working model is rewritten as:

$$(3.2) \quad \mathbf{Y} = \mathbf{X}_{n1} \boldsymbol{\beta}_1 + \mathbf{X}_{n2} \boldsymbol{\beta}_2 + \mathbf{X}_{n3} \boldsymbol{\beta}_3 + \boldsymbol{\varepsilon}, \quad p > n, \quad p_1 + p_2 < n.$$

**3.2. Overfitted Model.** We apply a variable selection method which keeps both strong and weak-moderate signals as follows:

$$(3.3) \quad \mathbf{Y} = \mathbf{X}_{n1} \boldsymbol{\beta}_1 + \mathbf{X}_{n2} \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}, \quad p_1 \leq p_2 < n.$$

Recall, the LASSO strategy which usually eliminates the sparse signals and retains weak-moderate and strong signals in the resulting model, and may be considered as an overfitted Model

**3.3. Underfitted Model.** Now, we apply a variable selection method which keeps only strong signals and eliminates all other signals in the resulting model. Thus, we have

$$(3.4) \quad \mathbf{Y} = \mathbf{X}_{n1} \boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}, \quad p_1 < n.$$

One can use ALASSO strategy which usually retain the strong signals and may produce a lower dimensional model as compared with LASSO. This model may bay termed as an underfitted Model.

We are interested in estimating  $\beta_1$  when  $\beta_2$  may be a null vector, but we are not sure. We suggest Stein-type shrinkage strategy for estimating  $\beta_1$  under this real situation. In essence we would like to combine estimates of the overfitted with the estimates of underfitted models to improve the efficiency of an underfitted model.

**3.4. Double Shrinking.** Ahmed and Yüzbaşı (2016a) defined a shrinkage estimator of  $\beta_1$  by combining overfitted model estimate  $\hat{\beta}_1^{\text{OF}}$  with the underfitted  $\hat{\beta}_1^{\text{UF}}$  as

$$(3.5) \quad \hat{\beta}_1^S = \hat{\beta}_1^{\text{UF}} + \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \left( 1 - (p_2 - 2)W_n^{-1} \right), \quad p_2 \geq 3,$$

where, the weight function  $W_n$  is defined by

$$W_n = \frac{n}{\hat{\sigma}^2} (\hat{\beta}_2^{\text{LSE}})' (\mathbf{X}_{S_2}' \mathbf{M}_1 \mathbf{X}_{S_2}) \hat{\beta}_2^{\text{LSE}},$$

and  $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_{S_1} (\mathbf{X}_{S_1}' \mathbf{X}_{S_1})^{-1} \mathbf{X}_{S_1}'$ ,  $\hat{\beta}_2^{\text{LSE}} = (\mathbf{X}_{S_2}' \mathbf{M}_1 \mathbf{X}_{S_2})^{-1} \mathbf{X}_{S_2}' \mathbf{M}_1 \mathbf{Y}$  and

$$\hat{\sigma}^2 = \frac{1}{n-1} (\mathbf{Y} - \mathbf{X}_{S_1} \hat{\beta}_1^{\text{UF}})' (\mathbf{Y} - \mathbf{X}_{S_1} \hat{\beta}_1^{\text{UF}}).$$

The  $\hat{\beta}_1^{\text{UF}}$  is the ALASSO estimator and  $\hat{\beta}_1^{\text{OF}}$  is the LASSO estimator.

Here, under the concept of double shrinking of Yüzbaşı and Arashi (2016), we define a family of double shrunken estimators

$$(3.6) \quad \begin{aligned} \hat{\beta}_1^{\text{FS}} &= \hat{\beta}_1^{\text{OF}} - \frac{(p_2 - 2)r(W_n)}{W_n} \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \\ &= \hat{\beta}_1^{\text{UF}} + \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \left( 1 - \frac{(p_2 - 2)r(W_n)}{W_n} \right), \quad p_2 \geq 3, \end{aligned}$$

where  $r(x)$  is a continuous, bounded and differentiable function of  $x$ .

For  $r(x) = 1$ , we get the result of Ahmed and Yüzbaşı (2016a).

In the spirit of Alam and Thompson (1969), we consider the function  $r(x) = 1/(1+x^{-1})$  to get

$$(3.7) \quad \hat{\beta}_1^{\text{FS1}} = \hat{\beta}_1^{\text{UF}} + \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \left( 1 - \frac{(p_2 - 2)}{1 + W_n} \right), \quad p_2 \geq 3,$$

Further, by the virtue of Gaussian kernel, we consider the function  $r(x) = \exp(-x^2)$  to get

$$(3.8) \quad \hat{\beta}_1^{\text{FS2}} = \hat{\beta}_1^{\text{UF}} + \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \left( 1 - \frac{(p_2 - 2) \exp(-W_n^2)}{W_n} \right), \quad p_2 \geq 3,$$

Lastly, we propose to use  $r(x) = \arctan(x)$ , which yields the following superior estimator

$$(3.9) \quad \hat{\beta}_1^{\text{FS3}} = \hat{\beta}_1^{\text{UF}} + \left( \hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}} \right) \left( 1 - \frac{(p_2 - 2) \arctan(W_n)}{W_n} \right), \quad p_2 \geq 3,$$

In the forthcoming section we will be analyzing the performance of  $\hat{\beta}_1^{\text{FS1}}$  and  $\hat{\beta}_1^{\text{FS3}}$  and compare with the superior estimator of Ahmed and Yüzbaşı (2016a), i.e., the PS estimator. In the conclusions, we will discuss about the usage of  $\hat{\beta}_1^{\text{FS2}}$ .

#### 4. THEORETICAL CONSIDERATIONS

In this section, we develop some properties of the proposed estimators. Because of the complexity, we only consider orthonormal design. Note that in general, the LASSO is not an oracle procedure and is not consistent, whereas the ALASSO has oracle properties. Honestly, our result, is restrictive.

For our purpose, we assume  $n^{-1}\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ . Under the specified partitioning,  $n^{-1}\mathbf{X}'_{S_i}\mathbf{X}_{S_i} = \mathbf{I}_{p_i}$  and  $\mathbf{X}'_{S_i}\mathbf{X}_{S_j} = \mathbf{0}$ , for  $i \neq j = 1, 2, 3$ . Simply,  $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_{S_1}\mathbf{X}'_{S_1}$ ,  $\mathbf{X}'_{S_2}\mathbf{M}_1\mathbf{X}_{S_2} = n\mathbf{I}_{p_2}$ , and

$$(4.1) \quad \hat{\beta}_2^{LSE} = n^{-1}\mathbf{X}'_{S_2}\mathbf{Y}, \quad W_n = \frac{1}{\hat{\sigma}^2}(\hat{\beta}_2^{LSE})'\hat{\beta}_2^{LSE}.$$

Further, we have

$$(4.2) \quad \begin{aligned} \hat{\beta}_1^{OF} &= \left( \text{sgn}(\hat{\beta}_j^{LSE}) \left( |\hat{\beta}_j^{LSE}| - \frac{\lambda}{2} \right)^+, j = 1, \dots, p_1 \right)' \\ \hat{\beta}_1^{UF} &= \left( \text{sgn}(\hat{\beta}_j^{LSE}) \left( |\hat{\beta}_j^{LSE}| - \frac{\lambda}{2|\hat{\beta}_j^{LSE}|} \right)^+, j = 1, \dots, p_1 \right)' \end{aligned}$$

For the true parameter value,  $|\hat{\beta}_{1j}| > \lambda/2$ , for  $j = 1, \dots, p_o < p_1$ , where  $p_o$  is the true parameter value for the active set  $\{j : \beta_j \neq 0, j = 1, \dots, p_1\}$ . Then, it is easy to see

$$\hat{\beta}_1^{OF} - \hat{\beta}_1^{UF} = \left( \frac{\lambda}{2} \text{sgn}(\hat{\beta}_j^{LSE}) \left( \frac{1}{|\hat{\beta}_j^{LSE}|} - 1 \right), j = 1, \dots, p_o \right)'.$$

Therefore, we can obtain the following bound

$$(4.3) \quad \mathbf{D} = \hat{\beta}_1^{OF} - \hat{\beta}_1^{UF} < \text{sgn}(\hat{\beta}_j^{LSE}) \left( 1 - \frac{\lambda}{2} \right)$$

if  $\lambda$  is suitably chosen such that  $\lambda < 2$ .

Define the risk function of any estimator  $\hat{\beta}_1$  of true parameter  $\beta_1$  by  $R(\beta_1; \hat{\beta}_1) = \mathbb{E}(\hat{\beta}_1 - \beta_1)'(\hat{\beta}_1 - \beta_1)$ . Note, here we have  $j = 1, \dots, p_o$ .

Then, under the orthonormal assumption, we have

$$(4.4) \quad \begin{aligned} R(\beta; \hat{\beta}_1^{FS}) - R(\beta; \hat{\beta}_1^{OF}) &= (p_2 - 2)^2 \mathbb{E} \left[ \frac{r^2(W_n)}{W_n^2} \mathbf{D}'\mathbf{D} \right] \\ &\quad - 2(p_2 - 2) \mathbb{E} \left[ \frac{r(W_n)}{W_n} (\hat{\beta}_1^{OF} - \beta_1)' \mathbf{D} \right] \\ &< \left( 1 - \frac{\lambda}{2} \right)^2 (p_2 - 2)^2 \mathbb{E} \left[ \frac{r^2(W_n)}{W_n^2} \sum_{j=1}^{p_o} \text{sgn}(\hat{\beta}_j^{LSE}) \right] \\ &\quad - 2(p_2 - 2) \left( 1 - \frac{\lambda}{2} \right) \mathbb{E} \left[ \frac{r(W_n)}{W_n} (\hat{\beta}_1^{OF} - \beta_1)' \text{sgn}(\hat{\beta}^{LSE}) \right] \end{aligned}$$

As  $n \rightarrow \infty$ ,  $\hat{\beta}_1^{OF} \xrightarrow{P} \beta_1$ . Hence, for sufficiently large samples size  $n$ ,  $R(\beta; \hat{\beta}_1^{FS}) < R(\beta; \hat{\beta}_1^{OF})$ , i.e., the proposed  $\hat{\beta}_1^{FS}$  outperforms  $\hat{\beta}_1^{OF}$  ( $\hat{\beta}_1^{FS} \succ \hat{\beta}_1^{OF}$ ) as soon as  $\sum_{j=1}^{p_o} \text{sgn}(\hat{\beta}_j^{LSE}) < 0$ , under a probabilistic sense. This scenario is independent of the choice of  $r(\cdot)$  and hence, all the shrinkage estimators outperform the over fitted model. Similar conclusion can be discovered for the under fitted model, with a slightly different condition.

In general,  $\hat{\beta}_1^{FS} \succ \hat{\beta}_1^{OF}$  iff for all  $r(\cdot)$ , we have

$$(4.5) \quad \mathbb{E} \left[ \frac{r(W_n)}{W_n} \left\{ (p_2 - 2) \frac{r(W_n)}{W_n} \mathbf{D} - 2(\hat{\beta}_1^{OF} - \beta_1) \right\}' \mathbf{D} \right] < 0$$

Let

$$(4.6) \quad \alpha = \frac{(\hat{\beta}_1^{OF} - \beta_1)' \mathbf{D}}{\mathbf{D}'\mathbf{D}}$$



Then,  $\alpha$  satisfies  $\sqrt{n}(\hat{\beta}_1^{\text{OF}} - \beta_1) = \alpha\sqrt{n}(\hat{\beta}_1^{\text{OF}} - \hat{\beta}_1^{\text{UF}}) = \alpha \left[ \sqrt{n}(\hat{\beta}_1^{\text{OF}} - \beta_1) - \sqrt{n}(\hat{\beta}_1^{\text{UF}} - \beta_1) \right]$ .

Let  $\lambda = o(\sqrt{n})$  and  $\lambda n^{(\gamma-1)/2} \rightarrow \infty$ . Using Theorem 2 of [Zou \(2006\)](#),  $\sqrt{n}(\hat{\beta}_1^{\text{UF}} - \beta_1) \xrightarrow{P} 0$ . Consequently  $\alpha \rightarrow 1$ . Now, we are ready to find the bound on  $r(W_n)/W_n$ .

Suppose  $r(\cdot) > 0$  and is concave. Then, using Lemma 1 of [Casella \(1990\)](#),  $r(W_n)/W_n$  is non-increasing. Hence, by (4.5),  $\hat{\beta}_1^{\text{FS}} \succ \hat{\beta}_1^{\text{OF}}$  if for every  $r$  function we have

$$(4.7) \quad \frac{r(W_n)}{W_n} < \frac{2}{p_2 - 2} \frac{(\hat{\beta}_1^{\text{OF}} - \beta_1)' \mathbf{D}}{\mathbf{D}' \mathbf{D}} \rightarrow \frac{2}{p_2 - 2}.$$

## 5. MONTE CARLO SIMULATION

We consider a Monte Carlo simulation, and simulate the response from the following model:

$$(5.1) \quad y_i = x_{1i}\beta_1 + x_{2i}\beta_2 + \dots + x_{pi}\beta_p + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

where  $\varepsilon_i$  are i.i.d.  $N(0, 1)$  and  $x_{ij} = (\xi_{(ij)}^1)^2 + \xi_{(ij)}^2$  with  $\xi_{(ij)}^1 \sim N(0, 1)$  and  $\xi_{(ij)}^2 \sim N(0, 1)$  for all  $i = 1, 2, \dots, n$ ,  $j = 1, 2, \dots, p$ .

We consider the regression coefficients are set  $\beta = (\beta'_1, \beta'_2, \beta'_3)' = (\Lambda'_{p_1}, \Delta'_{p_2}, \mathbf{0}'_{p_3})'$ , where,  $\Lambda_{p_1}$ ,  $\Delta_{p_2}$  and  $\mathbf{0}_{p_3}$  mean the vectors of  $\Lambda$ ,  $\Delta$  and 0 with dimensions  $p_1$ ,  $p_2$  and  $p_3$ , respectively. If  $\Delta = 0$ , then it indicates that the null hypothesis is true. On the other hand, the larger values  $\Delta$  indicate the degree of violation of null hypothesis.

In this simulation setting, we simulated 250 data sets consisting of  $n = 150$ ,  $\Lambda = 1, 2$ ,  $p_1 = 4$ ,  $p_2 = 4, 8, 16$  and  $p_3 = 200, 400, 800$ .

The performance of an estimator is evaluated by using relative mean squared error (RMSE) criterion. The RMSE of an estimator  $\beta_1^*$  with respect to  $\hat{\beta}_1^{\text{OF}}$  is defined as follows

$$(5.2) \quad \text{RMSE}(\beta_1^*) = \frac{\text{MSE}(\hat{\beta}_1^{\text{OF}})}{\text{MSE}(\beta_1^*)},$$

where  $\beta_1^*$  is one of the listed estimators. If the RMSE of an estimator is larger than one, it indicates that it is superior to  $\hat{\beta}_1^{\text{OF}}$ . The results of simulated RMSE of the listed estimators are reported in Tables 2 – 7 and Figures 1 and 2. We also report the TP (the number of true positives) and the FP (the number of false positives) in Table 1 only for  $(p_2, p_3) = (4, 200)$ .

According to the simulation results, the performance of under-fitted estimator ALASSO is the best since it is based on true model, and the FS3 performs better than both FS1 and PS when  $\Delta = 0$ . On the other hand, the RMSE of the ALASSO decreases and approaches to zero while the all others approach to one when we increase the magnitude of weak signals.

In Figure 1, if  $\Delta = 0$ , then both LASSO and ALASSO methods always select strong covariates, while ALASSO select less weak signals than LASSO. Contrary to this, if we increase the magnitude of weak signals, say  $\Delta = 0.8$ , then we observe that LASSO is more efficient than ALASSO for selecting those signals, see the Figure 2. For both case, our suggest methods again select all strong signals, while they select more variables than ALASSO when the weak signals are getting stronger.

Table 1 shows the numbers and the percentages of TP and FP of the listed estimators when  $(p_2, p_3) = (4, 200)$ . According to this table, all listed methods select all strong covariates for each values of  $\Delta$ , whereas ALASSO is the best for FP, which is indicated by the smallest ratio of FP.

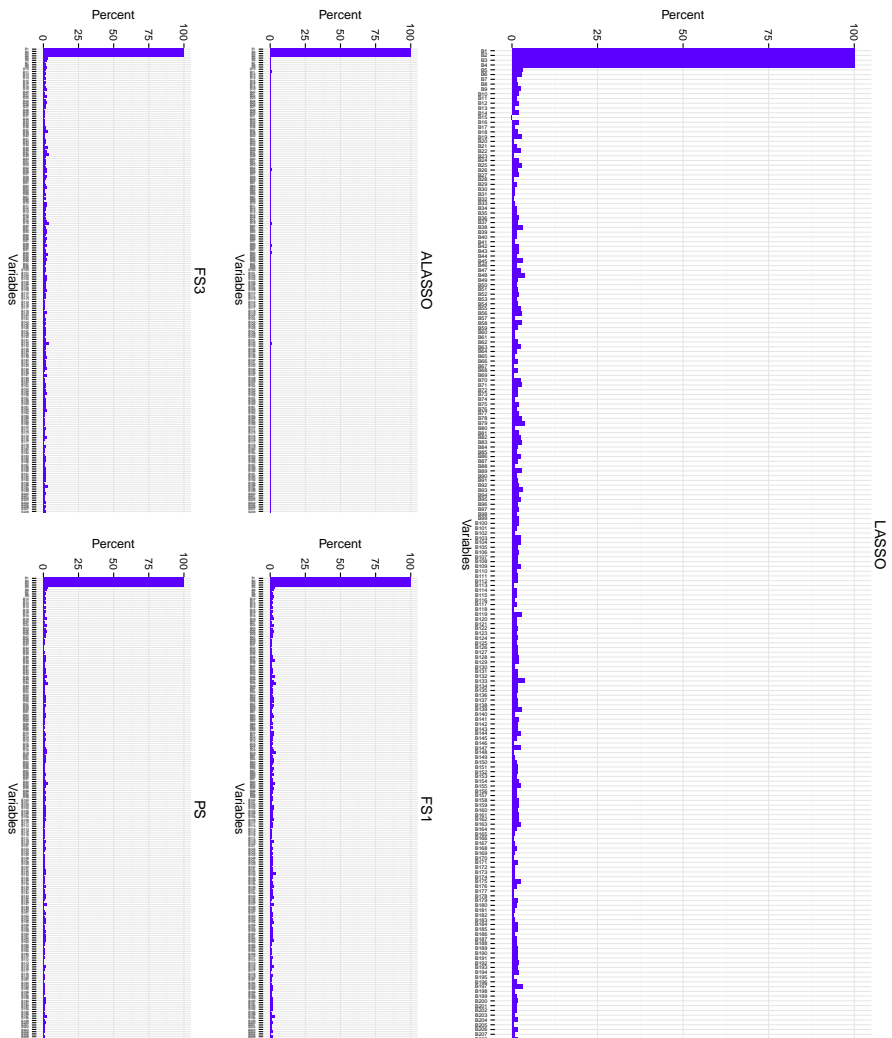


FIGURE 1. The percentage of times each predictor was selected when  $\Delta = 0$  and  $(p_2, p_3) = (4, 200)$

TABLE 1. The numbers and the percentages of TP and FP of estimators when  $(p_2, p_3) = (4, 200)$

$\Delta$	0.000	0.200	0.400	0.600	0.800	0.000	0.200	0.400	0.600	0.800
# the number of TP					# the number of FP					
LASSO	4.00	4.00	4.00	4.00	4.00	3.06	6.21	12.98	15.56	16.62
ALASSO	4.00	4.00	4.00	4.00	4.00	0.16	0.95	4.03	5.28	5.18
FS1	4.00	4.00	4.00	4.00	4.00	3.06	6.21	12.98	15.56	16.62
FS3	4.00	4.00	4.00	4.00	4.00	3.06	6.21	12.98	15.56	16.62
PS	4.00	4.00	4.00	4.00	4.00	2.06	6.21	12.98	15.56	16.62
# the percentages of TP					# the percentages of FP					
LASSO	100.00	100.00	100.00	100.00	100.00	0.76	1.54	3.21	3.85	4.11
ALASSO	100.00	100.00	100.00	100.00	100.00	0.04	0.24	1.00	1.31	1.28
FS1	100.00	100.00	100.00	100.00	100.00	0.76	1.54	3.21	3.85	4.11
FS3	100.00	100.00	100.00	100.00	100.00	0.76	1.54	3.21	3.85	4.11
PS	100.00	100.00	100.00	100.00	100.00	0.51	1.54	3.21	3.85	4.11

## 6. REAL DATA EXAMPLES

In this section, we analyze the performance of double shrinkage estimators dealing with two real data sets. We follow [Ahmed and Yüzbaşı \(2016a\)](#) strategies for analyzing each

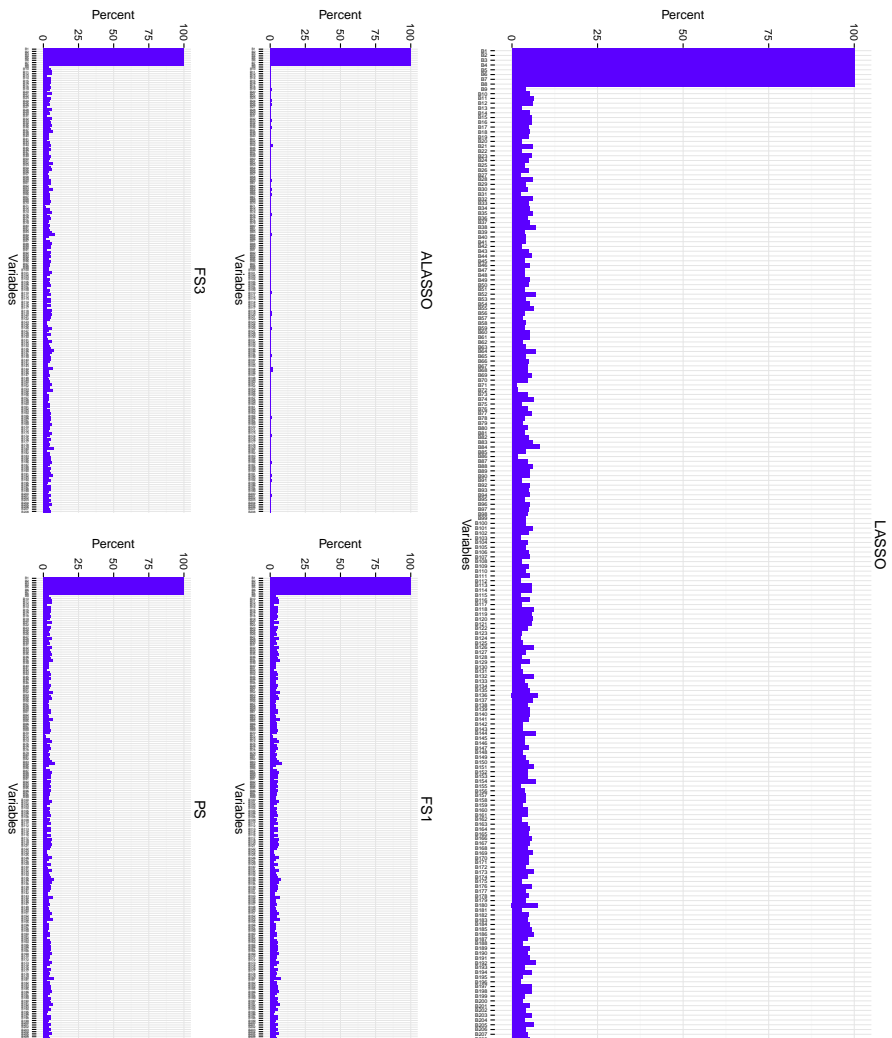


FIGURE 2. The percentage of times each predictor was selected when  $\Delta = 0.8$  and  $(p_2, p_3) = (4, 200)$

TABLE 2. RMSE of estimators when  $p_3 = 200$  and  $\Lambda = 1$

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	1.717	1.325	1.507	1.399
	0.200	1.324	1.040	1.065	1.043
	0.400	0.670	1.004	1.006	1.004
	0.600	0.324	0.994	0.991	0.994
8	0.800	0.189	0.992	0.988	0.992
	0.000	1.727	1.547	1.689	1.574
	0.200	1.238	1.051	1.080	1.052
	0.400	0.516	0.997	0.995	0.997
16	0.600	0.235	0.985	0.976	0.985
	0.800	0.142	0.984	0.974	0.984
16	0.000	1.749	1.677	1.888	1.653
	0.200	1.104	1.050	1.076	1.050
	0.400	0.458	0.993	0.988	0.993
	0.600	0.200	0.979	0.967	0.979
16	0.800	0.113	0.979	0.967	0.979

TABLE 3. RMSE of estimators when  $p_3 = 400$  and  $\Lambda = 1$ 

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	1.923	1.401	1.653	1.496
	0.200	1.483	1.046	1.073	1.048
	0.400	0.827	1.008	1.012	1.008
	0.600	0.382	0.996	0.994	0.996
	0.800	0.227	0.994	0.990	0.994
8	0.000	1.920	1.674	1.939	1.703
	0.200	1.332	1.061	1.097	1.063
	0.400	0.637	1.003	1.004	1.003
	0.600	0.295	0.989	0.983	0.989
	0.800	0.166	0.986	0.978	0.986
16	0.000	1.927	1.789	2.078	1.776
	0.200	1.213	1.056	1.087	1.057
	0.400	0.634	1.002	1.002	1.002
	0.600	0.287	0.988	0.980	0.988
	0.800	0.158	0.984	0.975	0.984

TABLE 4. RMSE of estimators when  $p_3 = 800$  and  $\Lambda = 1$ 

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	2.131	1.455	1.740	1.579
	0.200	1.587	1.052	1.083	1.055
	0.400	0.977	1.011	1.018	1.011
	0.600	0.468	0.998	0.997	0.998
	0.800	0.265	0.995	0.992	0.995
8	0.000	2.255	1.777	1.955	1.874
	0.200	1.472	1.069	1.110	1.071
	0.400	0.822	1.009	1.014	1.009
	0.600	0.382	0.994	0.990	0.994
	0.800	0.212	0.990	0.984	0.990
16	0.000	2.175	2.039	2.406	1.984
	0.200	1.325	1.065	1.103	1.066
	0.400	0.842	1.009	1.013	1.009
	0.600	0.468	0.996	0.994	0.996
	0.800	0.261	0.991	0.985	0.991

data sets. As long as there is no uncertain prior information about  $p_1$ ,  $p_2$  and  $p_3$ , one may use LASSO and ALASSO methods to find important covariates. After that, we may construct our estimation strategies. We also indicate that we draw 1000 bootstrap sample with dimension of the design matrix, and we calculate the prediction error (PE) based on 5 - fold cross validation, and take its average value for each bootstrap sample. To easy comparison, we report relative PE (RPE) of an estimator with respect to LASSO. Thus, a value of RPE  $> 1$  reflects the superiority of the other methods.

**6.1. Eye Data.** This data set contains gene expression data of mammalian eye tissue samples, [Scheetz et al. \(2006\)](#). The format is a list containing the design matrix which represents the data of  $n = 120$  rats with  $p = 200$  gene probes and the response vector with 120 dimensional which represents the expression level of TRIM32 gene. The numbers of selected variables for eye data set are 24 and 11 by LASSO and ALASSO, respectively.

TABLE 5. RMSE of estimators when  $p_3 = 200$  and  $\Lambda = 2$ 

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	1.799	1.356	1.596	1.436
	0.200	1.249	1.032	1.051	1.034
	0.400	0.690	1.003	1.005	1.003
	0.800	0.193	0.993	0.990	0.993
	1.200	0.085	0.993	0.990	0.993
	1.600	0.048	0.994	0.991	0.994
8	0.000	1.758	1.531	1.752	1.557
	0.200	1.081	1.031	1.048	1.032
	0.400	0.518	0.994	0.990	0.994
	0.800	0.134	0.984	0.974	0.984
	1.200	0.060	0.986	0.978	0.986
	1.600	0.034	0.988	0.981	0.988
16	0.000	1.766	1.689	2.048	1.652
	0.200	1.001	1.029	1.043	1.030
	0.400	0.472	0.992	0.987	0.992
	0.800	0.114	0.979	0.967	0.979
	1.200	0.051	0.982	0.972	0.982
	1.600	0.028	0.985	0.976	0.985

TABLE 6. RMSE of estimators when  $p_3 = 400$  and  $\Lambda = 2$ 

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	1.934	1.427	1.734	1.547
	0.200	1.459	1.044	1.070	1.047
	0.400	0.821	1.006	1.010	1.007
	0.800	0.224	0.994	0.991	0.994
	1.200	0.100	0.994	0.991	0.994
	1.600	0.060	0.995	0.992	0.995
8	0.000	2.052	1.743	2.041	1.784
	0.200	1.247	1.043	1.068	1.044
	0.400	0.639	1.002	1.002	1.002
	0.800	0.166	0.987	0.979	0.987
	1.200	0.073	0.988	0.981	0.988
	1.600	0.043	0.990	0.984	0.989
16	0.000	2.022	1.961	2.518	1.885
	0.200	1.080	1.034	1.052	1.035
	0.400	0.618	1.000	0.999	1.000
	0.800	0.159	0.984	0.975	0.984
	1.200	0.071	0.985	0.977	0.985
	1.600	0.040	0.987	0.980	0.987

In Figure 3, we plot the prediction error of each bootstrap replication for listed estimation techniques. Also, in Table 8, we report the RPEs of estimators. It can be seen that the performance of FS3 is the best which is followed by PS and FS1.

**6.2. Riboavin Data.** Here, we consider the data set about riboavin (vitamin B2) production in *Bacillus subtilis*. There is a single real valued response variable which is the logarithm of the riboavin production rate. Furthermore, there are  $p = 4088$  explanatory

TABLE 7. RMSE of estimators when  $p_3 = 800$  and  $\Lambda = 2$ 

$p_2$	$\Delta$	ALASSO	FS1	FS3	PS
4	0.000	2.396	1.524	1.903	1.686
	0.200	1.552	1.048	1.078	1.051
	0.400	0.943	1.009	1.014	1.009
	0.800	0.258	0.996	0.993	0.996
	1.200	0.121	0.995	0.992	0.995
	1.600	0.067	0.995	0.993	0.995
8	0.000	2.231	1.789	2.145	1.853
	0.200	1.361	1.058	1.092	1.059
	0.400	0.803	1.007	1.011	1.007
	0.800	0.217	0.990	0.984	0.990
	1.200	0.099	0.990	0.984	0.990
	1.600	0.053	0.991	0.986	0.991
16	0.000	2.203	2.093	2.677	2.000
	0.200	1.250	1.053	1.082	1.053
	0.400	0.875	1.009	1.014	1.009
	0.800	0.254	0.990	0.984	0.990
	1.200	0.107	0.989	0.983	0.989
	1.600	0.062	0.990	0.984	0.990

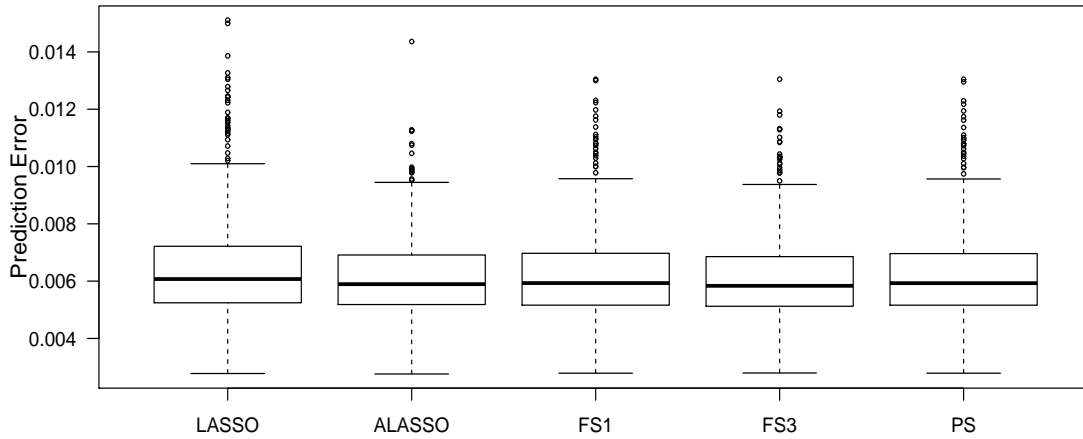


FIGURE 3. The RPEs for Eye data set.

TABLE 8. The average of RPEs for Eye data set.

ALASSO	FS1	FS3	PS
1.0598	1.0423	1.0626	1.0430

variables measuring the logarithm of the expression level of 4088 genes. There is one



rather homogeneous data set from  $n = 71$  samples that were hybridized repeatedly during a fed batch fermentation process where different engineered strains and strains grown under different fermentation conditions were analyzed.

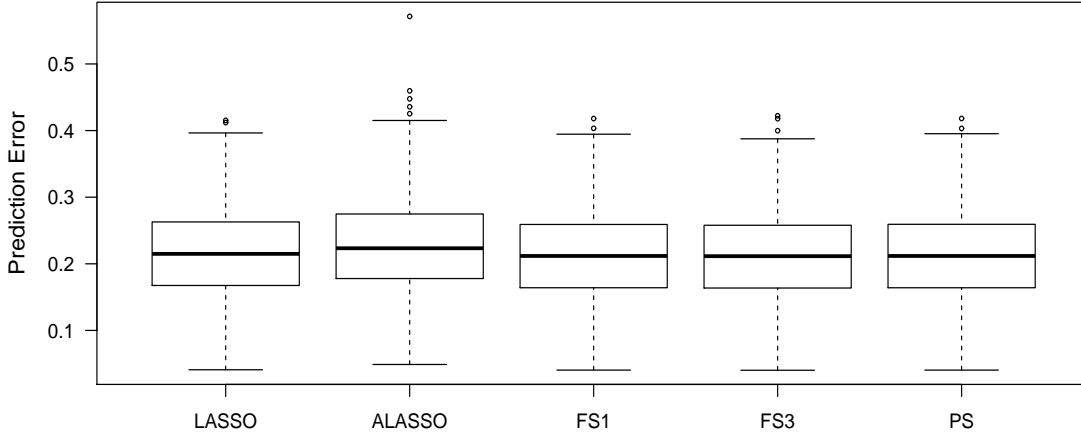


FIGURE 4. The average of RPEs for Riboflavin data set.

For this data, LASSO and ALASSO select 27 and 12 significant covariates, respectively. Notice that we used “one standard error” rule for selection of tuning parameters. We did not scale the design matrix and include the intercept term. The result of RPEs is shown in Table 4 and Figure 9. Again, it is clear that the performance of FS3 outshines PS and FS1 even though the performance of ALASSO is less efficient than LASSO.

TABLE 9. Relative Prediction Error of estimators

ALASSO	FS1	FS3	PS
0.9509	1.0222	1.0273	1.0225

## 7. CONCLUSIONS

In this paper, we extended variable selection methods to a new direction to be shrunken to a targeted estimator. Specifically we combined estimation strategies from both under-fitted and over-fitted models, in a high-dimensional regression model, employing a bounded measurable function. Specific concave functions were adopted to show the superiority of the proposed double shrunken estimators over the best of [Ahmed and Yüzbaşı \(2016a\)](#). We have conducted a simulation study to investigate the performance of the suggested shrinkage strategy with respect to two penalty estimators: LASSO and ALASSO. According to the simulation results, the performance of under-fitted estimator ALASSO is the best since it is based on true model, and the FS3 performs better than both FS1 and PS when  $\Delta = 0$ . On the other hand, the RMSE of the ALASSO decreases and approaches to zero while the all others approach to one when we increase the magnitude of weak signals. We further analysed two high-dimensional data sets, and the performance of the shrinkage strategy was striking.

We also proposed another shrinkage estimator, which was not included in the numerical analyses, namely FS2 given by

$$\widehat{\beta}_1^{\text{FS2}} = \widehat{\beta}_1^{\text{UF}} + \left( \widehat{\beta}_1^{\text{OF}} - \widehat{\beta}_1^{\text{UF}} \right) \left( 1 - \frac{(p_2 - 2) \exp(-W_n^2)}{W_n} \right), \quad p_2 \geq 3,$$

We evaluated the performance of FS2 comparatively. We realized, it is not competitor to FS1, FS3 and PS, specially for small values  $\Delta$ . However, as soon as the non-centrality parameter  $\Delta$  gets larger, its RMSE goes to 1 dominating all other estimators uniformly, since the weight goes to infinity and it simplifies to over-fitted model.

## REFERENCES

- Ahmed, S.E. (2014a). *Penalty, Shrinkage and Pretest Strategies: Variable Selection and Estimation*, Springer, New York.
- Ahmed, S.E. (2014b). *Perspectives on Big Data Analysis: Methodologies and Applications*, AMS, USA
- Ahmed, S. E., Hossain, S., & Doksum, K. A. (2012). LASSO and shrinkage estimation in Weibull censored regression models. *Journal of Statistical Planning and Inference*, 142(6), 1273-1284.
- Ahmed, S. E., & Yüzbaşı, B. (2016a). Big data analytics: integrating penalty strategies, *International Journal of Management Science and Engineering Management*, 11(2), 105-115
- Ahmed, S. E., & Yüzbaşı, B. (2016b). High Dimensional Data Analysis: Integrating Submodels. In *Big and Complex Data Analysis* (pp. 285-304). Springer International Publishing.
- Alam, K., & Thompson, J. R. (1969). Locally averaged risk, *Annals of Institute of Statistical Mathematics*, 21, 457-469.
- Armagan, A., Dunson, D. B., & Lee, J. (2013). Generalized double Pareto shrinkage. *Statistica Sinica*, 23(1), 119.
- Belloni, A., & Chernozhukov, V. (2009). Least squares after model selection in high-dimensional sparse models.
- Bhattacharya, A., Pati, D., Pillai, N. S., & Dunson, D. B. (2012). Bayesian shrinkage. arXiv preprint arXiv:1212.6088.
- Bickel, P.J., Ritov, Y. and Tsybakov, A.B. (2009). Simultaneous analysis of LASSO and dantzig selector, *Annals of Statistics*, 37, 1705–1732.
- Carvalho, C.M., Polson, N.G., & Scott, J.G. (2010). The horseshoe estimator for sparse signals, *Biometrika* 97, 465–480.
- Casella, G. (1990). Estimators with nondecreasing risk: Application of a chi-squared identity, *Statistics & Probability Letter* 10, 107–109.
- Efron, B., Hastie, T., Johnstone, I., & Tibshirani, R. (2004). Least angle regression. *The Annals of statistics*, 32(2), 407-499.
- Fan, J., & Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456), 1348-1360.
- Frank, I. and Friedman, J. (1993). A statistical view of some chemometrics regression tools (with discussion). *Technometrics*, 35, 109 – 148.
- Gao, X., Ahmed, S. E., & Feng, Y. (2016). Post Selection Shrinkage Estimation for High Dimensional Data Analysis. arXiv preprint arXiv:1603.07277.
- Hansen, B. E. (2015). The risk of James–Stein and Lasso shrinkage. *Econometric Reviews*, 1-15.
- Huang, J., Ma, S., & Zhang, C. H. (2008). Adaptive LASSO for sparse high-dimensional regression models. *Statistica Sinica*, 1603-1618.

- Kim, Y., Choi, H., & Oh, H. S. (2008). Smoothly clipped absolute deviation on high dimensions. *Journal of the American Statistical Association*, 103(484), 1665-1673.
- Leng, C., Lin, Y., & Wahba, G. (2006). A note on the LASSO and related procedures in model selection. *Statistica Sinica*, 1273-1284.
- Liu, H., & Yu, B. (2013). Asymptotic properties of lasso+ mls and lasso+ ridge in sparse high-dimensional linear regression. *Electronic Journal of Statistics*, 7, 3124-3169.
- Pyne, S., Rao, B.L.S. Prakasa and Rao, S.B.(2016). Big data analytics: Views from statistical and computational perspectives, in *Big Data Analytics: Methods and Applications*, 2016, Springer, India
- Saleh, A. K. Md. Ehsanes. (2006). *Theory of Preliminary Test and Stein-type Estimation with Applications*, John Wiley, New York.
- Schelldorfer, J., Bhlmann, P., DE, G., & VAN, S. (2011). Estimation for High Dimensional Linear Mixed Effects Models Using  $L_1$  - Penalization. *Scandinavian Journal of Statistics*, 38(2), 197-214.
- Scheetz, T. E., Kim, K. Y. A., Swiderski, R. E., Philp, A. R., Braun, T. A., Knudtson, K. L., ... & Sheffield, V. C. (2006). Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proceedings of the National Academy of Sciences*, 103(39), 14429-14434.
- Tibshirani, R. (1996). Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society. Series B (Methodological)*, 267-288.
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., & Knight, K. (2005). Sparsity and smoothness via the fused LASSO. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1), 91-108.
- Tran, M. N. (2011). The loss rank criterion for variable selection in linear regression analysis. *Scandinavian Journal of Statistics*, 38(3), 466-479.
- Wang, C., Chen, M.H., Schifano, E., Wu, J. and Yan, J. (2015). Statistical methods and computing for big data, arXiv:1502.07989.
- Wang, H., & Leng, C. (2012). Unified LASSO estimation by least squares approximation. *Journal of the American Statistical Association*.
- Yuan, M., & Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1), 49-67.
- Yuzbasi, B., & Arashi, M. (2016). Double shrunken selection operator, arXiv:1612.06304.
- Zhang, C.H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Annals os Statistics*, 38, 894 – 942.
- Zhang, C.H. and Zhang, S.S. (2014). Confidence intervals for low-dimensional parameters in high-dimensional linear models. *Annals os Statistics*, 76, 217–242.
- Zhao, P., & Yu, B. (2006). On model selection consistency of LASSO. *The Journal of Machine Learning Research*, 7, 2541-2563.
- Zou, H. (2006). The adaptive LASSO and its oracle properties. *Journal of the American statistical association*, 101(476), 1418-1429.